

CLAIMS

We claim:

1 1. A method for clustering a set S of n data points to find k final centers,
2 comprising:

3 partitioning said set S into P disjoint pieces S_1, \dots, S_P ;

4 for each said piece S_i , determining a set D_i of k intermediate centers;

5 assigning each data point in each piece S_i to the nearest one of said k intermediate
6 centers;

7 weighting each of said k intermediate centers in each set D_i by the number of
8 points in the corresponding piece S_i assigned to that center; and

9 clustering said weighted intermediate centers together to find said k final centers,
10 said clustering performed using a specific error metric and a clustering method A.

1 2. A method according to claim 1 further comprising:

2 merging said weighted centers into a single dataset D' prior to clustering.

1 3. A method according to claim 1 wherein P is sufficiently large enough such
2 that each piece S_i obeys the constraint $|S_i| < M$, where M is the size of a physical memory
3 or a portion thereof to be used in processing said each piece.

1 4. A method according to claim 1 wherein if P is not sufficiently large
2 enough such that each piece S_i obeys the constraint $|S_i| < M$, where M is the size of a
3 physical memory or a portion thereof to be used in processing said each piece, then
4 iteratively performing partitioning, determining, assigning, and weighting until the sets
5 D' of weighted intermediate centers generated thereby obeys the constraint $|D'| < M$.

1 5. A method according to claim 4 wherein said clustering is performed upon
2 iteratively obtained weighted intermediate clusters.

1 6. A method according to claim 4 wherein said set S is replaced by weighted
2 intermediate centers of the previous iteration when iteratively performing said
3 partitioning, determining, assigning, and weighting.

1 7. A method according to claim 1 wherein said determining is performed
2 using said specific error metric and said clustering method A.

1 8. A method according to claim 1 wherein said specific error metric is the
2 minimizing of the sum of the squares of the distances between points and their nearest
3 centers.

1 9. A method according to claim 1 wherein said specific error metric is the
2 minimizing of the sum of the distances between points and their nearest centers.

1 10. A method according to claim 1 wherein said clustering method is an
2 approximation-based method.

1 11. A method according to claim 8 wherein the distance is the Euclidean
2 distance.

1 12. A method according to claim 9 wherein the distance is the Euclidean
2 distance.

1 13. A method according to claim 1 further comprising:
2 considering a second set of data points for obtaining a second k final centers after
3 said set S is clustered;
4 repeating partitioning, determining, assigning and weighting for said second set of
5 data points; and
6 clustering weighted intermediate centers obtained from said second set of data
7 points together with said weighted intermediate centers obtained from said data set S, said
8 clustering performed using said specific error metric and said clustering method A.

1 14. A method according to claim 1 wherein said partitioning, determining,
2 assigning and weighting is performed in parallel for each piece S_i .

1 15. An article comprising a computer readable medium having instructions
2 stored thereon which when executed causes clustering a set S of n data points to find k
3 final centers, said clustering implemented by:

4 partitioning said set S into P disjoint pieces S_1, \dots, S_P ;

5 for each said piece S_i , determining a set D_i of k intermediate centers;

6 assigning each data point in each piece S_i to the nearest one of said k intermediate
7 centers;

8 weighting each of said k intermediate centers in each set D_i by the number of
9 points in the corresponding piece S_i assigned to that center; and

10 clustering said weighted intermediate centers together to find said k final centers,
11 said clustering performed using a specific error metric and a clustering method A.

1 16. An article according to claim 15 further implemented by:
2 merging said weighted centers into a single dataset D' prior to clustering.

1 17. An article according to claim 15 wherein P is sufficiently large enough
2 such that each piece S_i obeys the constraint $|S_i| < M$, where M is the size of a physical
3 memory or a portion thereof to be used in processing said each piece.

1 18. An article according to claim 15 wherein if P is not sufficiently large
2 enough such that each piece S_i obeys the constraint $|S_i| < M$, where M is the size of a
3 physical memory or a portion thereof to be used in processing said each piece, then
4 iteratively performing partitioning, determining, assigning, and weighting until the sets
5 D' of weighted intermediate centers generated thereby obeys the constraint $|D'| < M$.

1 19. An article according to claim 1 further implemented by:

2 considering a second set of data points for obtaining a second k final centers after
3 said set S is clustered;

4 repeating partitioning, determining, assigning and weighting for said second set of
5 data points; and

clustering weighted intermediate centers obtained from said second set of data

7 points together with said weighted intermediate centers obtained from said data set S, said
8 clustering performed using said specific error metric and said clustering method A,
9 resulting in said second k final clusters.

1 20. A method according to claim 1 wherein said partitioning, determining,
2 assigning and weighting is performed in parallel for each piece S_i .

1 21. An apparatus for clustering a set S of n data points to find k final centers,
2 said apparatus comprising:

3 a main memory;

4 a processor coupled to said memory, said processor configured to partition said set
5 S into P disjoint pieces S_1, \dots, S_P such that each piece S_i fits in main memory, said each
6 piece S_i first stored separately in said main memory and then clustered by said processor
7 performing:

8 for each said piece S_i , determining a set D_i of k intermediate centers;

9 assigning each data point in each piece S_i to the nearest one of said k intermediate
10 centers;

11 weighting each of said k intermediate centers in each set D_i by the number of
12 points in the corresponding piece S_i assigned to that center; and

13 clustering said weighted intermediate centers together to find said k final centers,
14 said clustering performed using a specific error metric and a clustering method A.

1 22. An apparatus for clustering a set S of n data points to find k final centers,

2. said apparatus comprising:

3 a main memory;

4 a plurality of processors coupled to said main memory, one of said processors
5 configured to partition said set S into P disjoint pieces S_1, \dots, S_P such that each piece S_i fits
6 in main memory, said each piece S_i first stored separately in said main memory and then
7 clustered by each said processor performing:
8 for each said piece S_i , determining a set D_i of k intermediate centers;
9 assigning each data point in each piece S_i to the nearest one of said k intermediate
10 centers; and
11 weighting each of said k intermediate centers in each set D_i by the number of
12 points in the corresponding piece S_i assigned to that center, further wherein after aid
13 weighting, one of said processors finally clustering said weighted intermediate centers
14 together to find said k final centers, said clustering performed using a specific error
15 metric and a clustering method A.